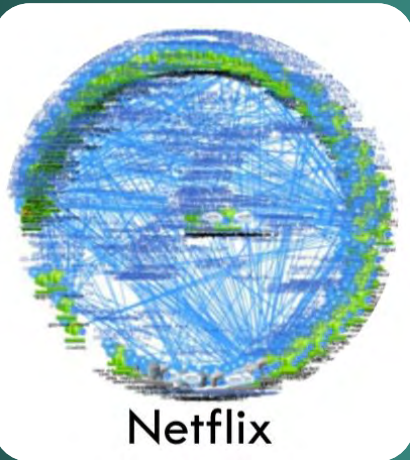# Aspects of Microservice Interactions

# Death Star Architecture



Netflix
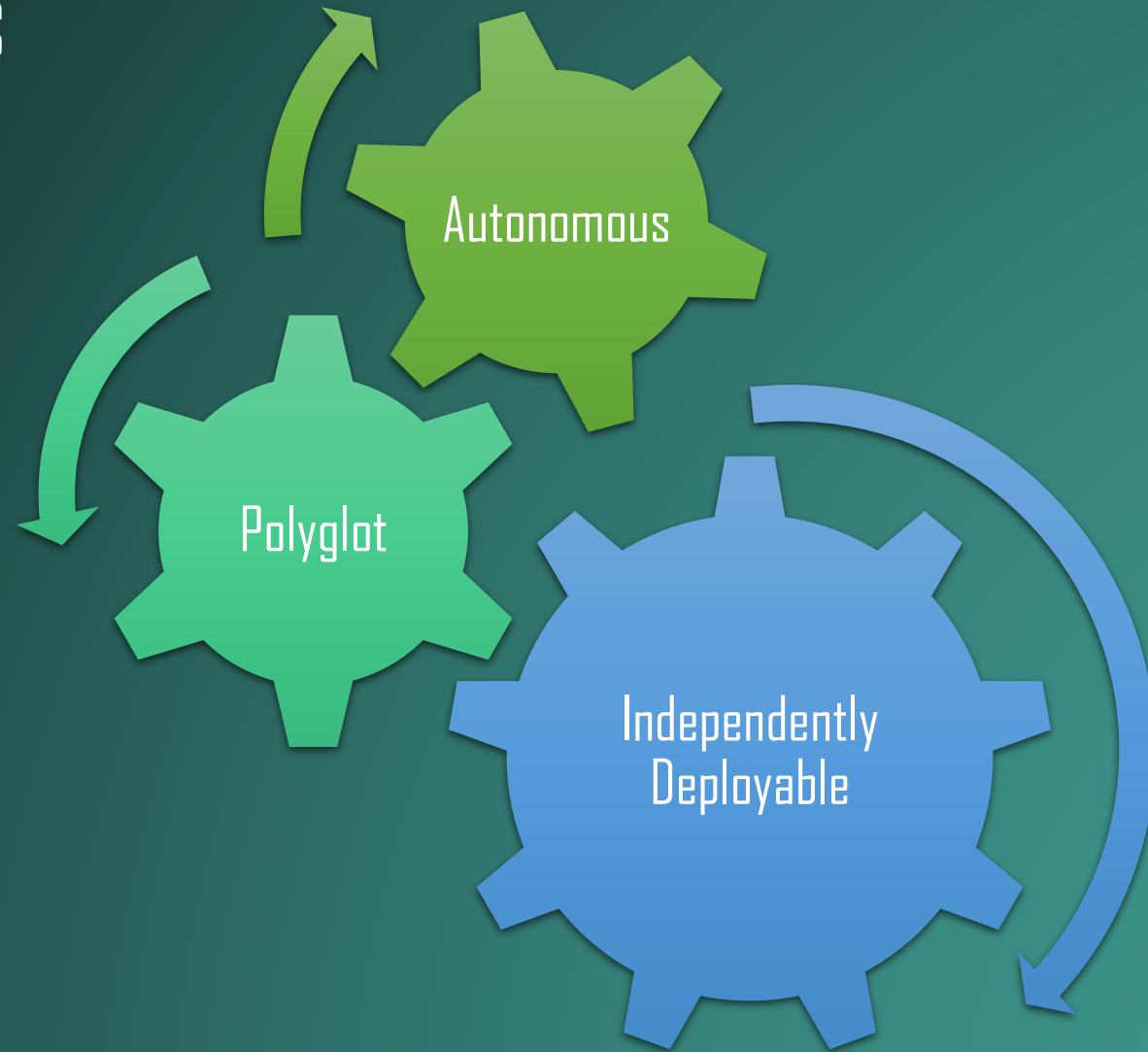


Twitter



Amazon

# WHY?
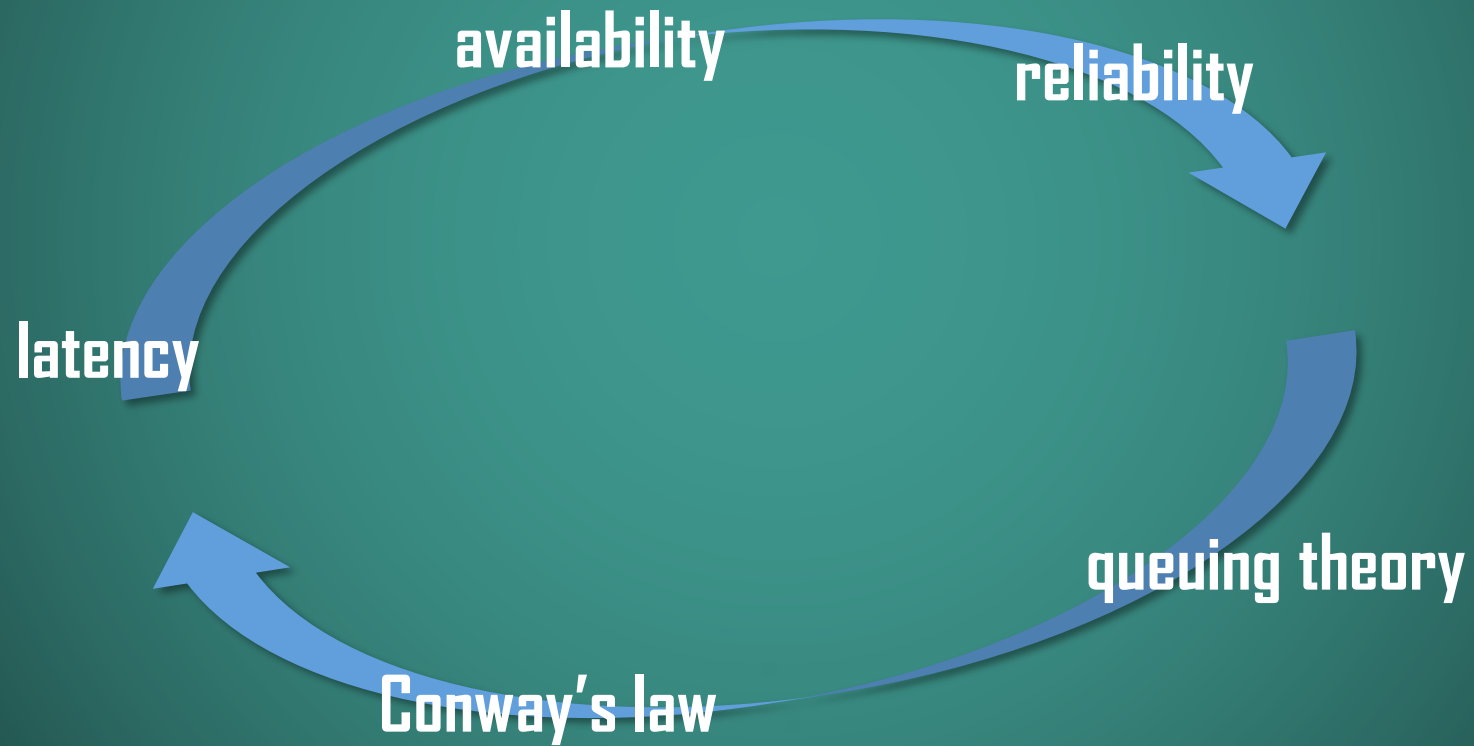
# Characteristics



Autonomous

Polyglot

Independently
Deployable

NETWORK CALL

# Driving Forces

availability

reliability

latency

queuing theory

Conway's law

# Latency



Radius: **4008**km / **2490.46**mi, RTT in vacuum: **27**ms, RTT in fiber: **41.04**ms.

# Latency

# Latency Numbers Every Programmer Should Know

Main memory reference: 100ns

1,000ns ≈ 1μs

Compress 1KB wth Zippy: 2,000ns ≈ 2μs

10,000ns ≈ 10μs =

Send 2,000 bytes over commodity network: 44ns

SSD random read: 16,000ns ≈ 16μs

Read 1,000,000 bytes sequentially from memory: 3,000ns ≈ 3μs

Round trip in same datacenter: 500,000ns ≈ 500μs

1,000,000ns = 1ms =

Read 1,000,000 bytes sequentially from SSD: 49,000ns ≈ 49μs

Disk seek: 2,000,000ns ≈ 2ms

Read 1,000,000 bytes sequentially from disk: 825,000ns ≈ 825μs

Packet roundtrip CA to Netherlands: 150,000,000ns ≈ 150ms

# Fallacies of distributed computing

文A  9 languages ∨

## The fallacies  [ edit ]

The fallacies are[1]

1. The network is reliable;
2. Latency is zero;
3. Bandwidth is infinite;
4. The network is secure;
5. Topology doesn't change;
6. There is one administrator;
7. Transport cost is zero;
8. The network is homogeneous.

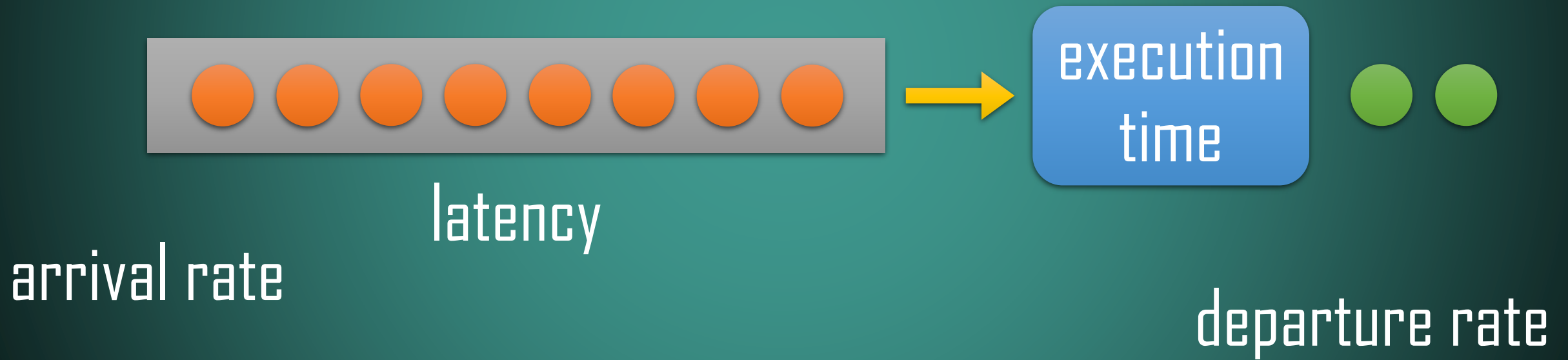# Queuing Theory

arrival rate → **backpressure!** → departure rate

# Queuing Theory

latency = 800ms

100ms

10/s

# Queuing Theory



latency = 800ms

50ms → 50ms → 20/s

# Queuing Theory

latency = 1200ms

100ms

50ms

10/s

# Queuing Theory

```
latency = sum(execution time / parallelism) * queue length



throughput = min((parallelism / execution time))
```

# Queuing Theory

QoS?

# Queuing Theory

Conway's Law

# Conway's Law

Conway's Law

# Conway's Law

Schema Owner

Conway's Law

Schema Federation

# TOOLBOX

CQRS
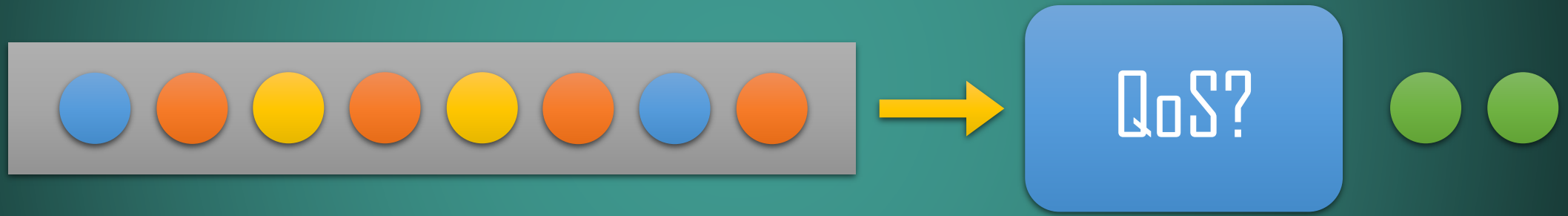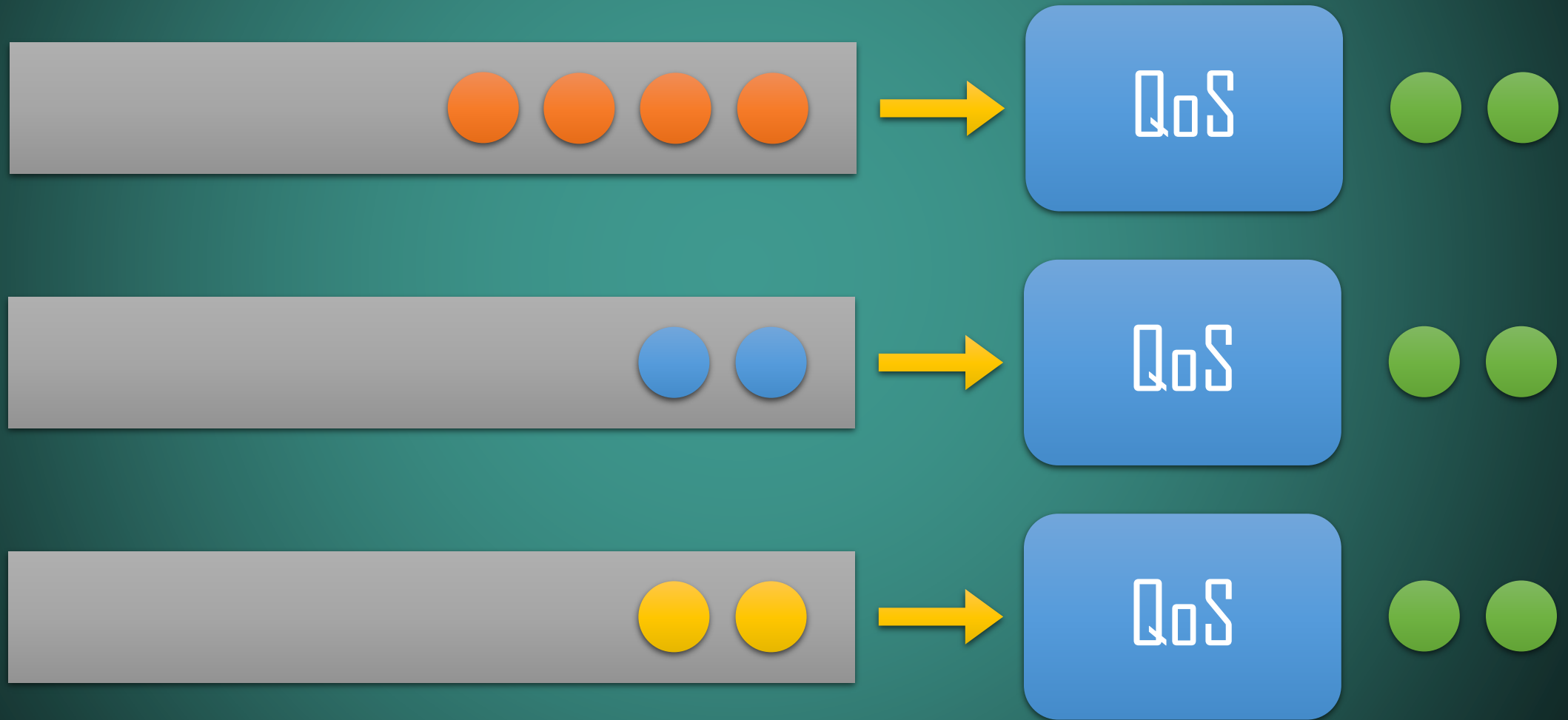
SCHEMAS
CONTRACT-BASED TESTS

CACHING
MULTI-WRITES

SYNC TO ASYNC
POLLING, FAKE BOUNDARIES

CLOUD NATIVE

AUTO-SCALING

BACKPRESSURE

SAGA
ORCHESTRATION, COREOGRAPHY

12 FACTOR

SERVICE MESH

TECHNOLOGY
GRPC, REST, GRAPHQL

MESSAGING
CORRELATION ID, ROUTING SLIP, EXACTLY ONCE
TOPICS VS QUEUES

CONCURRENCY
MODEL
REACTIVE, ACTOR-BASED, COROUTINES

RESILIENCY PATTERNS
CIRCUIT BREAKER, BULKHEAD, RETRIES, TIMEOUTS

OBSERVABILITY
SLO, TRACING, METRICS, LOGS

# EXAMPLES

Rapid Read Protection

Client

99%

99%

~ 99.99 %

‹epam›

- Read-heavy workload

- Immutable writes

- Low latency baseline

- Cold Cache & Distributed

  Cache won't work

READ    WRITE

‹epam›

# CQRS

WRITE  READ  READ  READ

Replication

Amazon Simple Queue
Service (Amazon SQS)

Amazon DynamoDB

hazelcast

‹epam›

# Client Libraries

Startup

**Amazon S3**

‹epam›

# Client Libraries

Startup

Amazon S3

# Client Libraries



Health Check

Amazon S3

# Shift Left

CI/CD

Amazon S3

‹epam›

# Proxy



Proxy

Amazon S3

Region Pinning

~60ms

Replication Lag

Region Pinning

~60ms

Replication Lag

Going...Going...Gone!

# Separate Critical Path



Creates an Executor that uses a single worker thread operating off an unbounded queue. (Note however that if this single thread terminates due to a failure during execution prior to shutdown, a new one will take its place if needed to execute subsequent tasks.) Tasks are guaranteed to execute sequentially, and no more than one task will be active at any given time. Unlike the otherwise equivalent newFixedThreadPool(1) the returned executor is guaranteed not to be reconfigurable to use additional threads.

Returns: the newly created single-threaded Executor

```
@NotNull
public static ExecutorService newSingleThreadExecutor() {
    return new FinalizableDelegatedExecutorService
        (new ThreadPoolExecutor( corePoolSize: 1,   maximumPoolSize: 1,
                        keepAliveTime: 0L, TimeUnit.MILLISECONDS,
                        new LinkedBlockingQueue<Runnable>()));
}
```

Recent Keys    Values

# Separate Critical Path

**Recent Keys**          **Values**

```
n unbounded queue. (
```

```
Creates a LinkedBlockingQueue with a capacity of Integer.MAX_VALUE.

public LinkedBlockingQueue() {
    this(Integer.MAX_VALUE);
}
```

# Baseline

**15,000 items** →

Auto Scaling group

HTTP 400

RCU = 400

# Retries

🚫 No Fairness

**Auto Scaling group**

Retry

HTTP 400

RCU = 400

# Built In Rate-Limiting



Auto Scaling group

40/s

40/s

80/s

HTTP 400

RCU = 400

‹epam›

# Built In Rate-Limiting

🚫 Rate-Limit Changes

Auto Scaling group

40/s

40/s

40/s

120/s

HTTP 400

RCU = 400

⟨epam⟩

# Service Mesh

Needs API Change

Istio

HTTP 429

120/s

Auto Scaling group

HTTP 400

RCU = 400

Queuing Theory

15,000 items

latency = 30s

4ms

4ms

500/s

‹epam›

# Moving Towards RabbitMQ

## Queue Overflow Behaviour

Use the `overflow` setting to configure queue overflow behaviour.

| | | | | |
|---|---|---|---|---|
| **persistent-q** | | D | | 🟨 flow | 117,215 |
| **transient-q** | | D | | 🟩 running | 0 |

```
Channel channel = ...;

Consumer consumer1 = ...;

Consumer consumer2 = ...;

channel.basicQos(10); // Per consumer limit

channel.basicConsume("my-queue1", false, consumer1);

channel.basicConsume("my-queue2", false, consumer2);
```

- `basic.ack` is used for positive acknowledgements
- `basic.nack` is used for negative acknowledgements (note: this is a RabbitMQ extension to AMQP 0-9-1)
- `basic.reject` is used for negative acknowledgements but has one limitation compared to `basic.nack`

# Single Worker

# Thank You!



**ORESZTÉSZ MARGARITISZ**

ASSOCIATE CHIEF SOFTWARE ENGINEER

@gitaroktato

gitaroktato

https://www.linkedin.com/in/oresztesz

# References

Latency
https://www.igvita.com/2012/07/19/latency-the-new-web-performance-bottleneck/
https://docs.google.com/a/chromium.org/viewer?a=v&pid=sites&srcid=Y2hyb21pdW0ub3JnfGRldnxneDoxMzcyOWIIN2I4YzI3NzE2
https://www.cloudping.co/
https://aws-latency-test.com/
https://colin-scott.github.io/personal_website/research/interactive_latency.html
https://gist.github.com/jboner/2841832
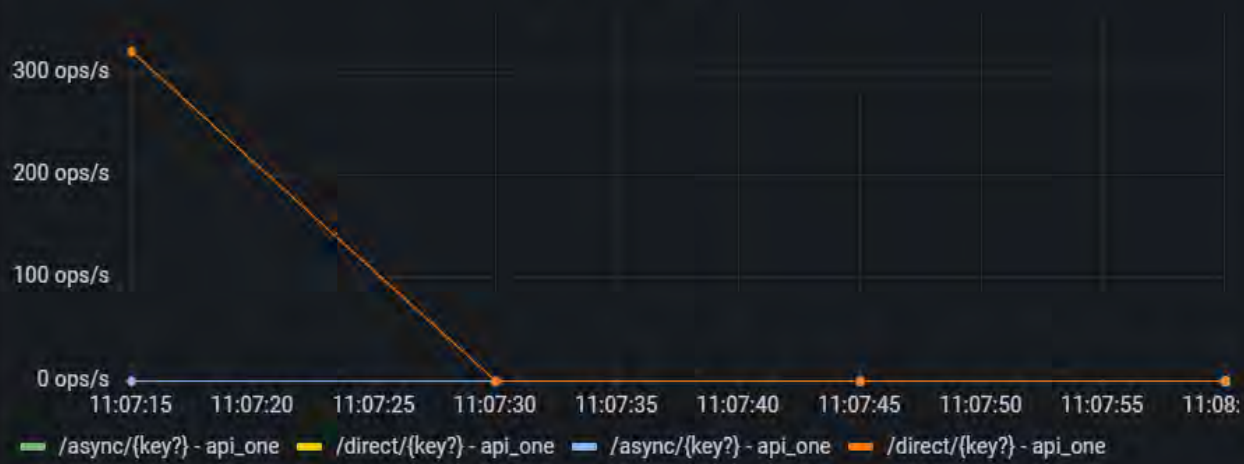https://blog.bytebytego.com/p/ep22-latency-numbers-you-should-know

Availability
https://github.com/gitaroktato/microservices-availability-simulator
https://eventhelix.com/fault-handling/system-reliability-availability
https://eventhelix.com/fault-handling/reliability-availability-basics

https://en.wikipedia.org/wiki/Fallacies_of_distributed_computing

# References

Backpressure with SQS
https://dl.awsstatic.com/architecture-diagrams/ArchitectureDiagrams/unreal-engines-pixel-streaming-on-aws-ra.pdf?did=wp_card&trk=wp_card

Reliability
https://www.joelonsoftware.com/2002/11/11/the-law-of-leaky-abstractions/

Queuing Theory
https://www.youtube.com/watch?v=oQGreeij-OE
https://www.youtube.com/watch?v=raRpbsWQBCo
https://en.wikipedia.org/wiki/Little%27s_law
https://en.wikipedia.org/wiki/Kendall%27s_notation
https://dzone.com/articles/applying-back-pressure-when

Rapid-Read Protection
http://www.datastax.com/dev/blog/rapid-read-protection-in-cassandra-2-0-2

# References

Architecture Katas
https://nealford.com/katas/kata?id=GoingGoingGone

Kafka / Kinesis Multi-Region Examples
https://aws.amazon.com/blogs/big-data/increase-apache-kafkas-resiliency-with-a-multi-region-deployment-and-mirrormaker-2
https://aws.amazon.com/blogs/big-data/build-highly-available-streams-with-amazon-kinesis-data-streams

The Twelve-Factor App
https://12factor.net/

Rate-Limiting Sandbox
https://github.com/gitaroktato/system-design-excercises/tree/main/rate-limiting

Microservices Availability Simulator
https://github.com/gitaroktato/microservices-availability-simulator

‹epam›

# References

Backpressure with RabbitMQ
https://www.rabbitmq.com/maxlength.html
https://www.rabbitmq.com/tutorials/tutorial-six-java.html
https://www.rabbitmq.com/consumers.html#single-active-consumer
https://www.rabbitmq.com/flow-control.html
https://blog.rabbitmq.com/posts/2020/05/quorum-queues-and-flow-control-the-concepts
https://www.rabbitmq.com/consumer-prefetch.html
https://blog.rabbitmq.com/posts/2014/04/finding-bottlenecks-with-rabbitmq-3-3/
https://blog.rabbitmq.com/posts/2015/10/new-credit-flow-settings-on-rabbitmq-3-5-5/
https://www.rabbitmq.com/confirms.html#publisher-confirms